

# The Role of Virtual Humans in Virtual Environment Technology and Interfaces

Daniel Thalmann  
Computer Graphics Lab, EPFL,  
Lausanne, Switzerland  
[thalmann@lig.di.epfl.ch](mailto:thalmann@lig.di.epfl.ch)  
<http://ligwww.epfl.ch>

## Abstract

The purpose of this paper is to show the importance of Virtual Humans in Virtual Reality and to identify the main problems to solve to create believable Virtual Humans.

## Introduction

The ultimate reason for developing realistic-looking Virtual Humans is to be able to use them in virtually any scene that re-creates the real world. However, a virtual scene -- beautiful though it may be -- is not complete without people.... Virtual people, that is. Scenes involving Virtual Humans imply many complex problems we have been solving for several years [1]. With the new developments of digital and interactive television [2] and multimedia products, there is also a need for systems that provide designers with the capability for embedding real-time simulated humans in games, multimedia titles and film animations. In fact, there are many current and potential applications of human activities that may be part of a VR system involving virtual humans:

- simulation based learning and training (transportation, civil engineering etc.)
- simulation of ergonomic work environments
- virtual patient for surgery, plastic surgery
- orthopedy and prostheses and rehabilitation
- plastic surgery
- virtual psychotherapies
- architectural simulation with people, buildings, landscapes and lights etc.
- computer games involving people and "Virtual Worlds" for Lunaparks/casinos
- game and sport simulation
- interactive drama titles in which the user can interact with simulated characters and hence be involved in a scenario rather than simply watching it.

But mainly, telepresence is the future of multimedia systems and will allow participants to share professional and private experiences, meetings, games, and parties. Virtual Humans have a key role to play in these shared Virtual Environments and true interaction with them is a great challenge. Although a lot of research has been going on in the field of Networked Virtual Environments, most of the existing systems still use simple embodiments for the representation of participants in the environments. More complex virtual human embodiment increases the natural interaction within the environment. The users' more natural perception of each other (and of autonomous actors) increases their sense of being together, and thus the overall sense of shared presence in the environment.



Fig.1 Virtual Humans

But, the modelling of Virtual Humans is an immense challenge as it requires to solve many problems in various areas. Table 1 shows the various aspects of research in Virtual Human Technology. Each aspects will be detailed and the problems to solve will be identified.

<b>Face and body representation</b>
<b>Avatar functions</b>
<b>Motion control</b>
<b>High-level behavior</b>
<b>Interaction with objects</b>
<b>Intercommunication</b>
<b>Interaction with user</b>
<b>Collaborative Virtual Environments</b>
<b>Crowds</b>
<b>Rendering</b>
<b>Standards</b>
<b>Applications</b>

Table 1. Aspects of research in Virtual Humans

### **Face and body representation**

Human modelling is the first step in creating Virtual Humans. For head, although it is possible to create them using an interactive sculpting tool, the best way is to reconstruct them from reality. Three methods have been used for this:

- 1) Reconstruction from 2D photos [3]
- 2) Reconstruction from a video sequence [4]
- 3) Construction based on the laser technology

The methods could be used for body modelling, but the main problem is still with the body deformations which has been addressed by many researchers, but is still not 100% solved.

Concerning facial expressions in Networked VEs, four methods are possible: video-texturing of the face, model-based coding of facial expressions, lip movement synthesis from speech and predefined expressions or animations. Believable facial emotions are still very hard to obtain.

*Main problem to solve:* realistic body and face construction and deformations

### **Avatar functions**

The avatar representation fulfils several important functions:

- 1) the visual embodiment of the user
- 2) means of interaction with the world
- 3) means of sensing various attributes of the world

It becomes even more important in multi-user Networked Virtual Environments [5], as participants' representation is used for communication. This avatar representation in NVEs has crucial functions in addition to those of single-user virtual environments [6 7]:

- 1) perception (to see if anyone is around)
- 2) localisation (to see where the other person is)
- 3) identification (to recognise the person)
- 4) visualisation of others' interest focus (to see where the person's attention is directed)
- 5) visualisation of other's actions (to see what the other person is doing and what is meant through gestures)
- 6) social representation of self through decoration of the avatar (to know what the other participants' task or status is)

Using articulated models for avatar representation fulfils these functionalities with realism, as it provides the direct relationship between how we control our avatar in the virtual world and how our avatar moves related to this control, allowing the user to use his/her real world experience. We chose to use complex virtual human models aiming for a high level of realism, but articulated "cartoon-like" characters could also be well suited to express ideas and feelings through the nonverbal channel in a more symbolic or metaphoric way.

*Main problem to solve:* easy way of directing an avatar

### **Motion control**

The main goal of computer animation is to synthesize the desired motion effect which is a mixing of natural phenomena, perception and imagination. The animator designs the object's dynamic behavior with his mental representation of causality. He/she imagines how it moves, gets out of shape or reacts when it is pushed, pressed, pulled, or twisted. So, the animation system has to provide the user with motion control tools able to translate his/her wishes from his/her own language.

In the context of Virtual Humans, a Motion Control Method (MCM) specifies how the Virtual Human is animated and may be characterized according to the type of information it privileged in animating this Virtual Human. For example, in a keyframe system for an articulated body, the privileged information to be manipulated is the angle. In a forward dynamics-based system, the privileged information is a set of forces and torques; of course, in solving the dynamic equations, joint angles are also obtained in this system, but we consider these as derived information. In fact, any MCM will eventually have to deal with geometric information (typically joint angles), but only geometric MCMs explicitly privilege this information at the level of animation control.

Many MCMs have been proposed: motion capture, keyframe, inverse kinematics, dynamics, walking models, grasping models, etc.. But, no method is perfect and only combination of blending of methods can provide good and flexible results.

*Main problem to solve:* flexible reuse, combination, and parameterisation of existing movements

### **High-level behavior**

Autonomous Virtual Humans should be able to have a behaviour, which means they must have a manner of conducting themselves. Typically, the Virtual Human should perceive the objects and the other Virtual Humans in the environment through virtual sensors [8]: visual, tactile and auditory sensors. Based on the perceived information, the actor's behavioural mechanism will determine the actions he will perform. An actor may simply evolve in his environment or he may interact with this environment or even communicate with other actors. In this latter case, we will consider the actor as a interactive perceptive actor.

Virtual vision [9] <sup>10</sup> is a main information channel between the environment and the virtual actor. The Virtual Human perceives his environment from a small window in which the environment is rendered from his point of view. As he can access depth values of the pixels, the colour of the pixels and his own position, he can locate visible objects in his 3D environment. To recreate the a virtual audition, in a first step, we have to model a sound environment where the Virtual Human can directly access to positional and semantic sound source information of a audible sound event. For virtual tactile sensors, it may be based on spherical multisensors attached to the articulated figure. A sensor is activated for any collision with other objects. These sensors could be integrated in a general methodology for automatic grasping.

A high level behavior uses in general sensorial input and special knowledge. A way of modeling behaviors is the use of an automata approach. Each actor has an internal state which can change each time step according to the currently active automata and its sensorial input. Abstraction mechanisms to simulate intelligent behaviours have been discussed in the AI (Artificial Intelligence) and AA (Autonomous Agents') literature. Several methods have been introduced to model learning processes, perceptions, actions, behaviours, etc, in order to build more intelligent and autonomous virtual agents.

*Main problem to solve:* development of very complex believable behaviors

## Interaction with objects

The necessity to model interactions between an object and a virtual human agent (here after just referred to as an agent), appears in most applications of computer animation and simulation. Such applications encompass several domains, as for example: virtual autonomous agents living and working in virtual environments, human factors analysis, training, education, virtual prototyping, and simulation-based design. A good overview of such areas is presented by Badler [11]. An example of an application using agent-object interactions is presented by Johnson et al [12], whose purpose is to train equipment usage in a populated virtual environment.

Commonly, simulation systems perform agent-object interactions for specific tasks. Such approach is simple and direct, but most of the time, the core of the system needs to be updated whenever one needs to consider another class of objects.

To overcome such difficulties, a natural way is to include within the object description, more useful information than only intrinsic object properties. Some proposed systems already use this kind of approach. In particular, the object specific reasoning [13] creates a relational table to inform object purpose and, for each object graspable site, the appropriate hand shape and grasp approach direction. This set of information may be sufficient to perform a grasping task, but more information is needed to perform different types of interactions.

Another interesting way is to model general agent-object interactions based on objects containing interaction information of various kinds: intrinsic object properties, information on how-to-interact with it, object behaviors, and also expected agent behaviors. The smart object approach, introduced by Kallmann and Thalmann [14 15] extends the idea of having a database of interaction information. For each object modeled, we include the functionality of its moving parts and detailed commands describing each desired interaction, by means of a dedicated script language. A feature modeling approach [16] is used to include all desired information in objects. A graphical interface program permits the user to interactively specify different features in the object, and save them as a script file.

*Main problem to solve:* make the Virtual Human learn how to interact with objects

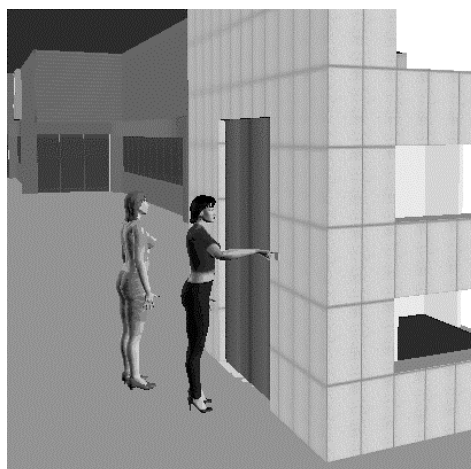


Fig.2. Interaction with objects

## Intercommunication

Behaviours may be also dependent on the emotional state of the actor. A non-verbal communication is concerned with postures and their indications on what people are feeling. Postures are the means to communicate and are defined by a specific position of the arms and legs and angles of the body. This non-verbal communication is essential to drive the interaction between people without contact or with contact.

What gives its real substance to face-to-face interaction in real life, beyond the speech, is the bodily activity of the interlocutors, the way they express their feelings or thoughts through the use of their body, facial expressions, tone of voice, etc. Some psychological researches have concluded that more than 65 percent of the information exchanged during a face-to-face interaction is expressed through nonverbal means [17]. A VR system that has the ambition to approach the fullness of real-world social interactions and to give to its participants the possibility to achieve a quality and realistic interpersonal communication has to address this point; and only realistic embodiment makes nonverbal communication possible.



Fig.3. Intercommunication

## Interaction with user

The real people are of course easily aware of the actions of the Virtual Humans through VR tools like Head-mounted displays, but one major problem to solve is to make the virtual actors conscious of the behaviour of the real people. Virtual actors should sense the participants through their virtual sensors. Such a perceptive actor would be independent of each VR representation and he could in the same manner communicate with participants and other perceptive actors. Perceptive actors and participants may easily be. For virtual audition, we encounter the same problem as in virtual vision. The real time constraints in VR demand fast reaction to sound signals and fast recognition of the semantic it carries. For the interaction between virtual humans and real ones, gesture recognition will be a key issue. As an example, Boulic et al. [18] produced a fighting between a real person and an autonomous actor. The motion of the real person is captured using a Flock of Birds. The gestures are recognised by the system and the information is transmitted to the virtual actor who is able to react to the gestures and decide which attitude to do.

*Main problem to solve:* development of a complete real-time vision-based recognition

## **Specific problems of Networked Virtual Environments**

Inserting virtual humans in the NVE is a complex task [5]. The main issues are:

- 1) selecting a scalable architecture to combine these two complex systems,
- 2) modeling the virtual human with believable appearance for interactive manipulation,
- 3) animating it with minimal number of sensors to have maximal behavioral realism,
- 4) investigating different methods to decrease the networking requirements for exchanging complex virtual human information.

Particularly, controlling the virtual human with limited input information is one of the main problems. For example, a person using a mouse will need extra input techniques or tools to exploit the functionalities of his embodiment. In this paper, we survey these tools that help a user with desktop VR configuration, we did not consider full tracking of the body using magnetic trackers, although this approach can be combined with limited tracking of the participant's arms.

*Main problems to solve:* controlling a realistic virtual human with limited input information, minimizing the information to be transmitted.

## **Crowds**

An accepted definition of crowd is that of a large group of individuals in the same physical environment, sharing a common goal (e.g. people going to a rock show or a football match). The individuals in a crowd may act in a different way than when they are alone or in a small group [19].

Although sociologists are often interested in crowd effects arising from social conflicts or social problems [20] the normal behavior of a crowd can also be studied when no changes are expected.

There are, however, some other group effects relevant to our work which are worth mentioning. Polarization occurs within a crowd when two or more groups adopt divergent attitudes, opinions or behavior and they may argue or fight even if they do not know each other. In some situations the crowd or a group within it may seek an adversary. The sharing effect is the result of influences by the acts of others at the individual level. Adding is the name given to the same effect when applied to the group. Domination happens when one or more leaders in a crowd influence the others.

Our goal [21] is to simulate the behavior of a collection of groups of autonomous virtual humans in a crowd. Each group has its general behavior [22] specified by the user, but the individual behaviors are created by a random process through the group behavior. This means that there is a trend shared by all individuals in the same group because they have a pre specified general behavior.

*Main problem to solve:* define collective behaviors while keeping individualities



Fig.4. Crowds

## Rendering

Rendering and animating in real-time a multitude of articulated characters presents a real challenge and few hardware systems are up to the task. Up to now little research has been conducted to tackle the issue of real-time rendering of numerous virtual humans. However, due to the growing interest in collaborative virtual environments the demand for numerous realistic avatars is becoming stronger.

There exist various techniques to speed up the rendering of a geometrical scene. They roughly fall into three categories: culling, geometric level-of-detail and image-based rendering which encompasses the concept of image caching. They all have in common the idea of reducing the complexity of the scene while retaining its visual characteristics.

Geometric level of detail (LOD) attempts to reduce the number of rendered polygons by using several representations of decreasing complexity of an object. At each frame the appropriate model or resolution is selected. Typically the selection criterion is the distance to the viewer although the object motion is also taken into account (motion LOD) in some cases. The major hindrance to using LOD is related to the problem of multi-resolution modeling, that is to say the automatic generation from a 3D object of simpler, coarser 3D representations that bear as strong a resemblance as possible to the original object.

Because 3D chips were not affordable or did not even exist in the 80s, video game characters, human-like or not, were then represented with 2D sprites. A sprite can be thought of as a block of pixels and a mask. The pixels give the color information of the final 2D image while the mask corresponds to a binary transparency channel. Using sprites, a human figure could easily be integrated into the decor. As more computing power was available in the 90s, the video game industry shifted towards 3D. However, the notion of sprites can also be used in the context of 3D rendering. This has been successfully demonstrated with billboards, which are basically 3D sprites, used for rendering very complex objects like trees or plants. In our opinion image-based rendering can also be used in the case of virtual humans by relying on the



intrinsic temporal coherence of the animation. Current graphics systems rarely take advantage of temporal coherence during animation. Yet, changes from frame to frame in a static scene are typically very small, which can obviously be exploited [23]. This still holds true for moving objects such as virtual humans providing that the motion remains slow in comparison with the graphics frame rate. Aubel and Thalmann [24] propose an approach to accelerated rendering of moving, articulated characters, which could easily be extended to any moving and/or self-deforming object. The method is based on impostors, a combination of traditional level-of-detail techniques and image-based rendering and relies on the principle of temporal coherence. It does not require special hardware (except texture mapping and Z-buffering capabilities, which are commonplace on high-end workstations nowadays) though fast texture paging and frame buffer texturing is desirable for optimal performance.



Fig.5. Use of Impostors

*Main problems to solve:* introduction of LOD for animation, integration with impostor technology

## **Standards**

Currently, a number of standardization efforts are continuing to solve different aspects of representing virtual humans in NVEs. Among them, the most significant efforts are MPEG-4 Face and Body Animation (FBA) [25], and VRML 2.0 Humanoid Animation (HAnim) [26] specifications. The VRML 2.0 HAnim group attempts to define a standard humanoid that can be exchanged between different users, or programs. The MPEG-4 FBA group aims at streaming virtual human bodies and faces with very low bitrates (less than 10 Kbits/second for the whole body).

*Main problem to solve:* defining standards for high-level behaviors of autonomous Virtual Humans

## **Areas of applications**

We may identify several areas [27] where autonomous virtual humans are essential:

*Virtual people for Inhabited Virtual Environments.* Their role is very important in virtual environments with many people, like virtual airports or even virtual cities. In the next few years, we will see a lot of Humanoids or Virtual Humans in many applications. These virtual humans will be more and more autonomous. They will also tend to become intelligent.

*Virtual substitutes.* A virtual substitute is an intelligent computer-generated agent able to act instead of the real person and on behalf of this person on the network. The virtual substitute has the voice of the real person and his or her appearance. He/she will appear on the screen of the workstation/TV, communicate with people, and have predefined behaviours planned by the owner to answer to the requests of the people.

*Virtual medical assistance.* Nowadays, it seems very difficult to imagine an effective solution for chronic care without including the remote care of patients at home by a kind of Virtual Medical Doctor. The modelling of virtual patient with correspondence to medical images is also a key issue and a basis for telesurgery.

## **Conclusions and recommendations**

Telepresence is the future of multimedia systems and will allow participants to share professional and private experiences, meetings, games, parties. The concepts of Distributed Virtual Environments are a key technology to implement this telepresence. Using humanoids within the shared environment is an essential supporting tool for presence. Real-time realistic 3D avatars will be essential in the future, but we will need interactive perceptive actors to populate the Virtual Worlds. The ultimate objective in creating realistic and believable virtual actors is to build intelligent autonomous virtual humans with adaptation, perception and memory. These actors should be able to act freely and emotionally. Ideally, they should be conscious and unpredictable. But, how far are we from such a ideal situation? Our interactive perceptive actors are able to perceive the virtual world, the people living in this world and in the real world. They may act based on their perception in an autonomous manner. Their intelligence is constrained and limited to the results obtained in the development of new methods of Artificial Intelligence. However, the representation under the form of virtual actors is a way of visually evaluating the progress. In the future, we may expect to meet intelligent actors able to learn or understand a few situations.

## **References**

1. N. Magnenat Thalmann, D. Thalmann, Complex Models for Animating Synthetic Actors, IEEE Computer Graphics and Applications, Vol.11, No5, 1991, pp.32-44.
2. N. Magnenat Thalmann N., Thalmann D. (1995) Digital Actors for Interactive Television, Proc. IEEE, Special Issue on Digital Television, Part 2, July 1995, pp.1022-1031.
3. W.S.Lee, N.Magnenat-Thalmann, Head Modeling from Pictures and Morphing in 3D with Image Metamorphosis Based on triangulation, in: Modelling and motion Capture Techniques for Virtual Environments, Lecture Notes in Artificial Intelligence, 1537, Springer, 1998.
4. P. Fua, R. Plankers and D. Thalmann, From Synthesis to Analysis: Fitting Human Animation Models to Image Data, Proc. CGI 99, IEEE Computer Society Press, 1999
5. T. K. Capin, I.S. Panzic, N. Magnenat-Thalmann, D. Thalmann Avatars in Networked Virtual Environments, John Wiley and Sons, 1999.

6. T. K. Capin, I.S. Panzic, N. Magnenat-Thalmann, D. Thalmann, Virtual Human Representation and Communication in VLNET Networked Virtual Environment, IEEE Computer Graphics and Applications, March 1997.
7. S. D. Benford et al., « Embodiments, Avatars, Clones and Agents for Multi-user, Multi-sensory Virtual Worlds », Multimedia Systems, Berlin, Germany: Springer-Verlag, 1997.
8. D. Thalmann, Virtual Sensors: A Key Tool for the Artificial Life of Virtual Actors, Proc. Pacific Graphics '95, Seoul, Korea, August 1995, pp.22-40.
9. O. Renault, N. Magnenat Thalmann, D. Thalmann (1990) A Vision-based Approach to Behavioural Animation, The Journal of Visualization and Computer Animation, Vol 1, No 1, pp 18-21.
10. X. Tu, D. Terzopoulos (1994) Artificial Fishes: Physics, Locomotion, Perception, Behavior, Proc. SIGGRAPH '94, Computer Graphics, pp.42-48.
11. N. N. Badler, "Virtual Humans for Animation, Ergonomics, and Simulation", IEEE Workshop on Non-Rigid and Articulated Motion, Puerto Rico, June 97.
12. W. L. Johnson, and J. Rickel, "Steve: An Animated Pedagogical Agent for Procedural Training in Virtual Environments", Sigart Bulletin, ACM Press, vol. 8, number 1-4, 16-21, 1997.
13. L. Levison, Connecting Planning and Acting via Object-Specific reasoning, PhD thesis, Dept. of Computer & Information Science, University of Pennsylvania, 1996.
14. M. Kallmann, D. Thalmann, Modeling Objects for Interaction Tasks, Proc. Eurographics Workshop on Animation and Simulation, Springer, 1998
15. M.Kallmann, D.Thalmann, A Behavioral Interface to Simulate Agent-Object Interactions in Real-Time, Proc. Computer Animation 99, IEEE Computer Society Press (to appear)
16. J.J. Shah, and M. Mäntylä, "Parametric and Feature-Based CAD/CAM", John Wiley & Sons, inc. 1995, ISBN 0-471-00214-3.
17. M. Argyle, Bodily Communication, New York: Methuen & Co., 1988.
18. L. Emering, R. Boulic, D. Thalmann, Interacting with Virtual Humans through Body Actions, IEEE Computer Graphics and Applications, 1998 , Vol.18, No1, pp.8-11.
19. M. E. Roloff. "Interpersonal Communication - The Social Exchange Approach". SAGE Publications, v.6, London. 1981.
20. J. S. McClelland. The Crowd and The Mob. Book printed in Great Britain at the University Press, Cambridge. 1989.
21. S.R. Musse, C. Babski, T. Capin, D. Thalmann, Crowd Modelling in Collaborative Virtual Environments, ACM VRST /98, Taiwan
22. S.R. Musse, D. Thalmann, A Model of Human Crowd Behavior: Group Inter-Relationship and Collision Detection Analysis. Proc. Workshop of Computer Animation and Simulation of Eurographics'97, Sept, 1997. Budapest, Hungary.
23. G. Schaufler, W. Stürzlinger, A Three Dimensional Image Cache for virtual reality, Proc. Eurographics'96, pp. C-227 – C-234
24. A.Aubel, R. Boulic, D.Thalmann, Animated Impostors for Real-time Display of Numerous Virtual Humans, Proc. Virtual Worlds 98, Paris, September 1998.
25. MPEG-N1902, Text for CD 14496-2 Video, ISO/IEC JTC1/SC29/WG11 N1886, MPEG97/November 1997.
26. Hanim, <http://ece.uwaterloo.ca/~h-anim/spec.html>
27. D.Thalmann, L.Chiariglione, F.Fluckiger, E.H. Mamdani, M.Morganti, J.Ostermann, J.Sesena, L.Stenger, A.Stienstra, Report on Panel 6: From Multimedia to Telepresence,

Expert groups in Visionary Research in Advanced Communications, ACTS, European Commission, 1997.